

# Self-supervised Monocular Image Depth Learning and Confidence Estimation

Long Chena, Wen Tanga, Tao Ruan Wan<sup>b</sup>, Nigel W. John<sup>c</sup>

<sup>a</sup> Bournemouth University, Poole, UK

<sup>b</sup> University of Bradford, Bradford, UK

<sup>c</sup> University of Chester, Chester, UK

---

## Abstract

We present a novel self-supervised framework for monocular image depth learning and confidence estimation. Our framework reduces the amount of ground truth annotation data required for training Convolutional Neural Networks (CNNs), which is often a challenging problem for the fast deployment of CNNs in many computer vision tasks. Our DepthNet adopts a novel fully differential patch-based cost function through the Zero-Mean Normalized Cross Correlation (ZNCC) to take multi-scale patches as matching and learning strategies. This approach greatly increases the accuracy and robustness of the depth learning. Whilst the proposed patch-based cost function naturally provides a 0-to-1 confidence, it is then used to self-supervise the training of a parallel network for confidence map learning and estimation by exploiting the fact that ZNCC is a normalized measure of similarity which can be approximated as the confidence of the depth estimation. Therefore, the proposed corresponding confidence map learning and estimation operate in a self-supervised manner and is a parallel network to the DepthNet. Evaluation on the KITTI depth prediction evaluation dataset and Make3D dataset show that our method outperforms the state-of-the-art results.

---

## Keywords:

Monocular depth estimation

Deep convolutional neural networks

Confidence map

## 1. Introduction

The human vision system is amazingly complex and extremely delicate. It can perceive depth through stereopsis, which relies on the displacement of the same object between the images received by the left and right retinas [1]. With extensive visual experience and through trial and error, humans develop the ability to use contextual depth cues to achieve good and reliable perception of depth and better understanding of spatial structure. Among these depth cues, most of them do not rely on stereopsis (the perception of depth from binocular vision), such as object occlusion, perspective, familiar and relative size, depth from motion, lighting and shading. Therefore, if blind in one eye or if performing a monocular task such as endoscopic surgery, we can still judge distance from these many different intuitive depth cues. In contrast, when using machine vision it is hard to infer the non-stereopsis depth cues.

With the recent development of Deep Convolutional Neural Networks (DCNNs), machines can solve many computer vision problems when provided with very large human annotated

datasets such as ImageNet [2], which is known as supervised learning. Acquisition of labelled datasets is one of the biggest challenges for supervised learning, however, which is an expensive, time-consuming and labour-intensive task.

In this paper, we propose a novel self-supervised computational framework that mimics the process of how a human learns various of contextual depth cues from stereopsis. We propose to “teach” the neural networks to “learn” the depth by themselves from “looking” to stereo image pairs. To be more specific, we construct a patch-based loss function that leverages the epipolar constraint [3] of stereo vision to minimize the depth prediction error from the context of a single image for each training iteration. Our approach does not require the ground truth depth for supervised training. Instead, we derive the implicit function of estimating depth from monocular images by the epipolar constraint of the stereo image pair, which is very easy to acquire compared with the ground truth depth that can only be obtained from LiDAR measurements. Therefore, our method can be regarded as self-supervised learning.

Compared with previous work [4–6] addressing the same problem, we propose a novel patch-based depth learning strategy, inspired by the classic patch matching algorithms for finding the best-matched patches between the left and right images. We

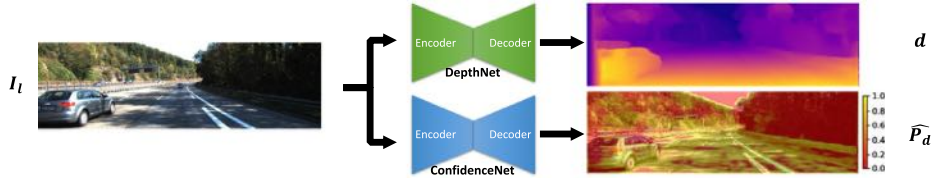


Fig. 1. Our proposed framework can simultaneously estimate depth and the confidence of estimated depth.

use the Zero-Mean Normalized Cross Correlation (ZNCC) to measure the normalized similarities between these patches. A fully-differential patch-based ZNCC cost function is implemented to guide the depth synthesis process for more accurate and robust results. Visual assessment shows that our approach can produce more accurate and reliable depth estimations in both texture-rich and texture-less areas due to the enlargement of matching field from a pixel to a patch (see Fig. 5). Empirical evaluations on KITTI dataset demonstrate the effectiveness of our approach and produce a state-of-the-art performance in monocular depth estimation task.

Our second contribution is that we train a parallel DCNN to evaluate the performance of the monocular depth estimation which can output a 0 to 1 confidence map. The parallel DCNN is also trained in a self-supervised manner thanks to our ZNCC similarity measurement function. As ZNCC is a normalized measure of similarity, which can be approximated as the confidence of the depth estimation, we take the ZNCC loss to self-supervise the parallel DCNN (ConfidenceNet) during training so that we can estimate the confidence of the depth estimated from the first DCNN (DepthNet) during testing mode as shown in Fig. 1. A confidence map is extremely useful for the monocular depth estimation task trained in an unsupervised manner, as the learned epipolar constraint only works well when there are clear corresponding pixels between the image pairs; it will fail and produce uncertain depth when occlusion and specularities exist in the images. Our confidence map can give a real-time assessment of the reliability of the predicted depth, which can then be further integrated into many applications such as monocular dense reconstruction [7], SLAM-based depth fusion [8], and many tasks need crucial accurate and confidence such as the monocular endoscopic surgery and the perception task for self-driving.

## 2. Related work

### 2.1. Stereo depth estimation

The problem of stereo images depth estimation has been well studied for a long time [9,10]. With the theory of epipolar constraint, accessing depth from stereo images can be regarded as a well-posed problem when ignoring the occlusions and depth discontinuities. Many stereo vision algorithms managed to achieve comparable results to ground truth depth acquired from depth sensors [11,12].

### 2.2. Monocular depth estimation

In contrast, estimating depth from monocular images is an ill-posed problem that is inherently ambiguous [13], and many research efforts have been devoted to the problem of monocular image depth estimation. One of the classic methods is Shape from Shading (SFS) [14], which is based on the gradual variation of shading as a cue to estimate the shape and depth. However, SFS has a strict prior assumption of Lambertian reflectance, uniform color and texture, and fixed light source direction, which are not applicable to most of the images in the real world. Saxena et al.

[15–18] used Markov Random Field (MRF) incorporated with multi-scale image features to learn monocular cues in a supervised manner. However, the hand-craft local features used in these approaches limit the expressive power of supervised learning, and lack a global contextual understanding of the scene for learning consistent depth.

### 2.3. DCNNs based monocular depth learning

More recently, DCNNs [13,19] are introduced to solve the challenge of monocular depth estimation problem, and has pushed the state-of-the-art forward in this area. Building on the success of this approach, several improvements have been made by incorporating probabilistic models such as Conditional Random Fields (CRFs) [20–24], advanced network structures such as Resnet [25], fully convolutional Resnet [26], two-streamed networks [27], multi-task joint training [19,28–31] and novel loss functions such as sparse semi-supervision [32,33], relative depth [34,35] and depth as classification [26]. Impressive as these works are, ground-truth depth data are still needed for the supervision of training these DCNNs. Recently,

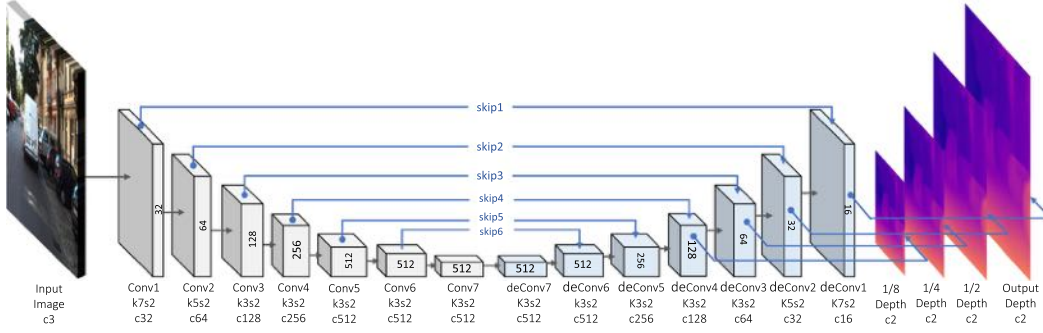
### 2.4. Unsupervised monocular depth learning

Driven by DCNNs, view synthesis technology [36] has proven to be effective on synthesizing new views by sampling pixels from existing views [37,38], which enables novel frameworks of unsupervised learning of monocular depth from stereo pairs, e.g., Deep3D [39], Garg et al. [4]. The works by Godard et al. [5] and Zhou et al. [6] advanced the networks by incorporating left-right consistency and pose estimations. Further improvements including introducing Visual Odometry (VO) or Multi-View Stereo (MVS) to learn depth from monocular videos [40–43]. However, a common weakness of these approaches is the use of pixel-wised photometric loss (L1-norm) to construct loss functions to guide the back-propagation process. Gradients are derived from the pixel intensity difference [6], which will lead to ambiguous gradients in texture-less areas and also in the regions that contain the mixture of thin structures and texture-less areas. Although multi-scale and smoothness loss functions are used to prevent such issue [4–6], the results are still not desirable and gradients are still likely to converge to local minimums due to the ambiguous pixel-wise loss. As shown in Fig. 5, in a common speed limitation board area from the KITTI dataset, the direct pixel-wise photometric loss will lead to many local minimums shown in the right curve chart. While as the left curve chart shows the result of using our proposed patch-based ZNCC loss, the loss is more smooth and likely to converge to the global minimum in the epipolar line. And the experiment result (the last row in Fig. 5) shows our proposed method can effectively generate accurate depth in complex regions.

### 2.5. Novelty compared to previous work

We propose a novel multi-scale patch-based cost function that adopts the ZNCC as a similarity function to explicitly enlarge the





**Fig. 3.** Depth synthesis network structure. “k” is the kernel size, “s” for the stride, “c” for the channel number. For simplicity, we do not draw the conv layers after each conv and deconv layer, which have the same kernel and channel size as previous layers but with stride 1.

in stereo vision. According to the epipolar constraint: the projection of a pixel  $x_l$  on the right camera plane  $x_r$  must be contained in the epipolar line. For calibrated stereo pairs discussed in this paper,  $x_l$  and  $x_r$  must be in the same row  $y$ , and the disparity  $d$  describes the horizontal displacement of the corresponding pixels  $x_l$  and  $x_r$ . Through the stereo triangulation, we can get that

$$D_{xy} = \frac{bf}{d} \Rightarrow d = x_l - x_r = \frac{bf}{D_{xy}} \quad (2)$$

where  $D_{xy}$  is the depth estimated in the pixel at  $(x, y)$ ,  $b$  and  $f$  are the camera baseline and focal distance. By the relationship discussed in the above equation, the target view in a stereo pair can be reconstructed given the source view and the corresponding depth (estimated through our depth synthesis network).

However, the direct mapping from one known view to the other view (forward mapping) will result in holes in the target image that are not differentiable. Therefore, we use the inverse mapping: for each pixel in the target view, by picking points from the source to reconstruct the target view guided by the  $d$ . Thus, a complete and differentiable target view can be generated. Then the bilinear sampling [49] is used to get the interpolated pixel value from the source view.

### 3.4. Disparity-guided patch sampling

Inspired by the stereo view reconstruction described above, we propose a novel patch sampling process guided by the estimated disparity from our DepthNet.  $N_{x,y}$  is defined as a patch with window size  $n$ , centered at the coordinate  $(x, y)$ . We sample patches on each pixel in the left image  $\{x, y \in I_l | I_l(N_{x,y})\}$ , and the corresponding patches shifted by disparity values  $d$  of each pixel in the right image,  $\{x, y \in I_r | I_r(N_{x-d,y})\}$ . According to Eq. 2, if  $d$  is correct, then we have  $I_l(N_{x,y}) = I_r(N_{x-d,y})$ . And this relationship will be used to construct the patch matching loss. These sampled patches are computed and stored vectorized so that can be deployed parallelly on GPU for accelerated computation.

The patch sampling size is very important and can affect the final performance of similarity measurement. However, there is no optimal patch size and the performance varies greatly across different images and local details. When small patch size is used, little information will be captured, and the similarity comparison robustness will be decreased. If we use a large patch size, computational complexity will be greatly increased and also cannot recover accurate depth at stereo occlusion and depth discontinuous. Therefore, we use a multi-scale patch sampling scheme and sample a combination of 4 different patch sizes in an image to fully exploit the effects of different patch sizes. We will discuss the choice of patch sizes in Section 4.1.3.

### 3.5. Loss function construction

We define a loss function  $L_{total}$  with multiple strategies to effectively train our networks for accurate, smooth and realistic depth.

$$L_{total} = \omega_p L_{PM} + \omega_v L_{VR} + \omega_d L_{DS} + \omega_c L_{DC} \quad (3)$$

where from left to right is: Patch Matching Loss, View Reconstruction Loss, Disparity Smoothness Loss and Disparity Consistency Loss.  $\omega$  is the corresponding weights to balance the effects of gradients back propagation. Each loss function will be explained in details below:

#### 3.5.1. Patch matching loss

Inspired by patch matching algorithm that by finding the best-matched patches in the left and right image to get correct disparities. We propose a patch matching loss that maximize the similarities (minimize the differences) of patches in left image  $I_l(N_{x,y})$  and the shifted patches in right image  $I_r(N_{x-d,y})$  to get correct disparities. Here, the ZNCC measure of similarity is used to compute a normalized similarity between the patches  $I_l(N_{x,y})$  and  $I_r(N_{x-d,y})$ :

$$C_{ZNCC}(I_l(N_{x,y}), I_r(N_{x-d,y})) = \frac{\sum_{i,j \in N_{x,y}} (I_l(i,j) - \bar{I}_l(N_{x,y})) \cdot (I_r(i-d,j) - \bar{I}_r(N_{x-d,y}))}{\sqrt{\sum_{i,j \in N_{x,y}} (I_l(i,j) - \bar{I}_l(N_{x,y}))^2 \cdot \sum_{i,j \in N_{x,y}} (I_r(i-d,j) - \bar{I}_r(N_{x-d,y}))^2}} \quad (4)$$

where  $\bar{I}(N_{x,y}) = \frac{1}{n} \sum_{x,y \in N_{x,y}} I(x,y)$  is the mean intensity of the patch  $N_{x,y}$  centered at the coordinate  $(x, y)$ .

The ZNCC returns a similarity ranging from  $[-1, 1]$ . We first normalize it into  $[0, 1]$  then invert it to get the patch matching loss:

$$L_{PM} = \sum_{x,y} 1 - \frac{1 + C_{ZNCC}(I_l(N_{x,y}), I_r(N_{x-d,y}))}{2} \quad (5)$$

Our patch matching loss is computed at all 4 patch sizes to cover both small structures and large areas. There are several advantages of using our patch-based ZNCC loss to regularize the depth synthesis:

(1) Our patch matching loss uses patches for measurement that involve larger regions than the direct pixel-wise photometric loss used in previous work, which is more robust and can achieve sub-pixel accuracy. Fig. 5 demonstrates the effect of our patch-based ZNCC loss. We charted the values of our patch-based ZNCC loss and the photometric loss against the disparity value of a pixel located at the center of the image patch “6”. It is obvious that by using our proposed patch-based ZNCC loss, the loss is more smooth

and likely to converge to the global minimum. Whereas the direct pixel-wise photometric loss will lead to many local minimums shown in the right curve chart.

(2) Compared to other similarity measures such as Sum of Absolute Differences (SAD), Census, and Normalized Cross Correlation (NCC), ZNCC is especially robust against Gaussian noise and variation between the compared patches, which can help to recover more accurate depth in our self-supervised framework.

(3) As a zero-mean normalized similarity measurement function, our patch-based ZNCC loss can provide a similar value ranging from  $[-1, 1]$ . After normalized to  $[0, 1]$  as shown in Eq. (5), it can be regarded as the confidence of the generated depth at each pixel, which can be further used to self-supervise the training of our confidence network.

### 3.5.2. View reconstruction loss

We use the view reconstruction loss as a second supervision on the depth synthesis. Guided by the synthesized depth, the right views can be reconstructed by collecting pixels from left images. The view reconstruction loss is defined as the L1 loss between the reconstructed view  $\hat{I}_r$  and the original view  $I_r$ :

$$L_{VR} = \sum_{xy} |I_r(x, y) - \hat{I}_r(x, y)| \quad (6)$$

Compared to the patch matching loss, the view reconstruction L1 loss is more sensitive to small structures and depth discontinuities and can provide more detailed depth information.

### 3.5.3. Disparity smoothness loss

We use a disparity smoothness term to regularize our network to produce more smooth depth. Similar to [4–6], we use the sum of the L1 norm of the disparity gradients along the  $x$  and  $y$  directions as a smoothness factor. The edge-aware terms are used to reduce the penalty on edges where depth discontinuities usually happen, which can prevent over-smoothing.

$$L_{DS} = \frac{1}{XY} \sum_{x,y} \left| \frac{\partial d(x, y)}{\partial x} \right| e^{-\left\| \frac{\partial d(x, y)}{\partial x} \right\|} + \left| \frac{\partial d(x, y)}{\partial y} \right| e^{-\left\| \frac{\partial d(x, y)}{\partial y} \right\|} \quad (7)$$

### 3.5.4. Disparity consistency loss

The left-right disparity consistency loss proposed in [5] has achieved a great improvement for monocular depth generation. Here, we adopt this loss function into our framework. The left and right image disparities are both generated, and the difference of left disparity map and the reconstructed left disparity map from right disparity is computed and minimized. This loss will ensure the left and right disparities coherence.

$$L_{DC} = \frac{1}{XY} \sum_{x,y} |d_l(x, y) - d_r(x - d_l(x, y), y)| \quad (8)$$

### 3.6. Confidence estimation network

One of the advantages of our proposed patch matching loss is that a normalized similarity measurement can be generated for each pixel at the training time. With the well-known epipolar constraint, the per-pixel confidence of the estimated depth can be approximated as the normalized similarity measurement of the left patches and the corresponding patches in the right image.

$$P_d(x, y) \approx C_{\text{Normalized}}(I_l(N_{x,y}), I_r(N_{x-d,y})) = (1 - L_{PM}(x, y)) \quad (9)$$

Here, we propose to use another encoder-decoder network to learn the confidence map generated by our depth estimation network during training, so that the confidence map can be preserved and generated during the testing time. We tried to train the confidence and depth in one network like [19, 28–30], but the multi-task

training would reduce the depth estimation performance. Therefore, we use a parallel encoder-decoder network to learn the confidence supervised by the per-pixel ZNCC loss of our depth estimation network. The loss of our ConfidenceNet is shown below:

$$L_{\text{ConfidenceNet}} = \sum_{x,y} |(1 - L_{PM}(x, y)) - \hat{P}_d(x, y)| \quad (10)$$

where  $\hat{P}_d(x, y)$  is the generated confidence map,  $L_{PM}(x, y)$  is the patch matching loss from our depth estimation network described in above sections. The static copy is used here to prevent the gradients propagating back to the depth estimation network. The  $1 - L_{PM}(x, y)$  operation inverts the loss to confidence, and L1 loss is used to access the confidence estimation error.

Instead of using the same encoder-decoder network structure as our DepthNet, we employ a simpler structure by only using first 5 conv-layer and last 5 deconv-layer without skip layers as described in Fig. 3 for two reasons:

(1) To reduce memory usage and training time, as training two neural networks at the same time is very computationally expensive. The second network can be replaced by a deeper and more complex encoder-decoder network to produce sharper and more accurate confidence, but the main purpose of our work is to prove that our self-supervised monocular depth learning and confidence estimation framework is feasible and helpful for depth prediction, hence we choose to use a simple network structure as the proof of concept.

(2) We intend to use a simpler network with fewer weights to prevent over-fitting to noises and to learn more generic confidence – high confidence in texture-rich areas, low confidence in texture-less, blurry and occluded areas, which is what we design this confidence net for.

## 4. Experiments

In this section, we evaluate our framework and compare the results with prior approaches both quantitatively and qualitatively on KITTI dataset. We use the rectified stereo image pairs for training our networks. For testing time, we use the left image to generate depth, and the corresponding sparse LIDAR data is served as the ground truth for benchmarking.

### 4.1. Implementation details

Our networks are implemented in Tensorflow and trained on a workstation with a single Nvidia Titan X GPU (12G Memory). Our models take around 60 hours to train for 50 epochs. When in testing mode, our networks can output depth and confidence map at around 20 frames per second.

#### 4.1.1. Hyper parameters

All input images are scaled to  $512 \times 256$  with a batch size of 4. Adam Optimizer is used with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and initial learning rate  $\lambda = 0.0001$  that decays after half of the training process. The weights to construct our total loss function for depth estimation network are  $w_p = 0.5$ ,  $w_v = 1$ ,  $w_d = 0.1$ ,  $w_c = 1$ .

#### 4.1.2. Data augmentation

The same data augmentation approach in [5] is used to randomly flip the image and change the gamma, brightness, and color shifts to increase the network robustness and prevent over-fitting.

#### 4.1.3. Multi-scale implementation

We employ a multi-scale strategy to ensure a coarse-to-fine up-sampling. As can be seen from Figs. 3, 4 depth scales are outputted with 1/8, 1/4, 1/2 and a full resolution. All of our loss functions are computed for each of these 4 scales, and for each of left and right



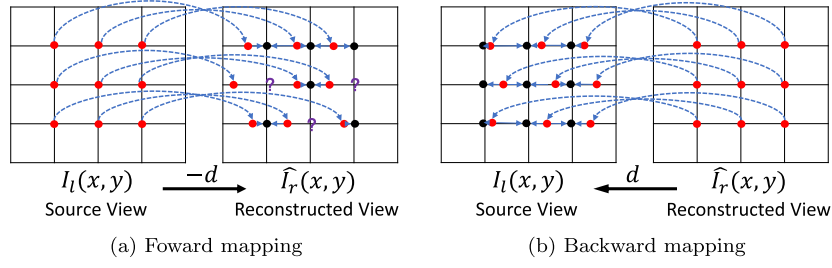


Fig. 4. The difference between forward mapping and backward mapping.

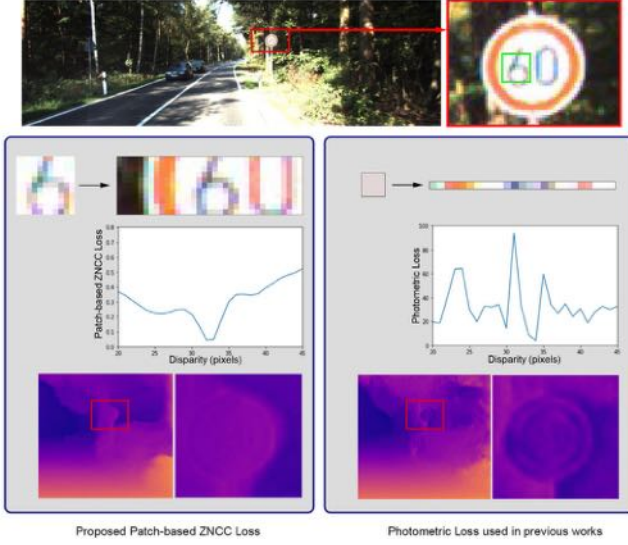


Fig. 5. Comparison of our proposed patch-based ZNCC loss (left image) with the photometric loss used in previous works (right image) to demonstrate that a patch naturally encodes more information than a single pixel and our loss function is more smooth and convex than other methods, therefore is more likely to converge to global minimum in the epipolar line.

images/disparities. We take the means of these loss functions as the final loss.

#### 4.1.4. Patch size

By applying different patch sizes on different image scales, we can get very large equivalent patch sizes with less computation. For patch size choices, based on our empirical test, we use

$n = 5, 5, 7, 9$  pixels for our patch-based ZNCC loss on 4 different scales, which is equivalent  $n = 5, 10, 28, 72$  pixels' windows on full resolution images.

#### 4.2. Training dataset

To be able to compare with the state-of-the-art monocular depth learning approaches, we trained and evaluated our networks using two different train/test splits: *Godard* and *Eigen*.

##### 4.2.1. Godard split

We use the same train/test sets that Godard et al. [5] proposed in their work. 200 high quality disparity images in 28 scenes provided by the official KITTI training set are served as the ground truth for benchmarking. For the rest of 33 scenes with a total of 30,159 images, 29,000 images are picked for training and the remaining 1159 images for testing.

##### 4.2.2. Eigen split

For fair comparison with more previous works, we also use the test split proposed by Eigen et al. [13] that has been widely evaluated by the works of Garg et al. [4], Liu et al. [23], Zhou et al. [6] and Godard et al. [5]. This test split contains 697 images of 29 scenes. The rest of 32 scenes contain 23,488 images, in which 22,600 are used for training and the remaining for testing, similar to [4] and [5].

#### 4.3. Results

##### 4.3.1. Quantitative evaluation

**Evaluation Metrics.** To access the quantitative performance of our proposed depth prediction network and compare with previous works, we evaluate each method using several error and accuracy metrics from [4–6,13]. The error metrics we use include

Table 1  
Comparison with state-of-the-art methods on KITTI dataset.

Method	Super-vision	Split	Cap	Error (Lower better)					Accuracy (Higher better)		
				AbsRel	SqRel	RMSE	RMSElog	D1-all	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Eigen et al. [13]	Yes	E	80	0.203	1.548	6.307	0.282	-	0.702	0.890	0.958
Liu et al. [23]	Yes	E	80	0.201	1.584	6.471	0.273	-	0.680	0.898	0.967
Zhou et al. [6]	No	E	80	0.208	1.768	6.856	0.283	-	0.678	0.885	0.957
Godard et al. [5]	No	E	80	0.148	1.344	5.927	0.247	-	0.803	0.922	0.964
ours-SAD	No	E	80	0.147	1.302	5.901	0.246	-	0.805	0.922	0.964
<b>ours-ZNCC</b>	<b>No</b>	<b>E</b>	<b>80</b>	<b>0.145</b>	<b>1.267</b>	<b>5.786</b>	<b>0.244</b>	-	<b>0.811</b>	<b>0.925</b>	<b>0.965</b>
Garg et al. [4]	No	E	50	0.169	1.080	5.104	0.273	-	0.740	0.904	0.962
Zhou et al. [6]	No	E	50	0.201	1.391	5.181	0.264	-	0.696	0.900	0.966
Godard et al. [5]	No	E	50	0.140	0.976	4.471	0.232	-	0.818	0.931	0.969
ours-SAD	No	E	50	0.140	0.959	4.463	0.232	-	0.821	0.931	0.969
<b>ours-ZNCC</b>	<b>No</b>	<b>E</b>	<b>50</b>	<b>0.138</b>	<b>0.937</b>	<b>4.399</b>	<b>0.231</b>	-	<b>0.825</b>	<b>0.933</b>	<b>0.969</b>
Godard et al. [5]	No	G	80	0.124	1.388	6.125	0.217	30.272	0.841	0.936	0.975
ours-SAD	No	G	80	0.121	1.358	6.073	0.215	29.937	0.842	0.936	0.976
<b>ours-ZNCC</b>	<b>No</b>	<b>G</b>	<b>80</b>	<b>0.117</b>	<b>1.202</b>	<b>5.953</b>	<b>0.210</b>	<b>29.612</b>	<b>0.845</b>	<b>0.938</b>	<b>0.976</b>

**Table 2**  
Comparison with state-of-the-art methods on Make3D dataset [18].

Method	Supervision	Cap	Error (Lower better)			
			AbsRel	SqRel	RMSE	RMSElog
Karsch et al. [50]	Yes	70	0.428	5.079	8.389	0.149
Liu et al. [21]	Yes	70	0.475	6.562	10.05	0.165
Laina et al. [25]	Yes	70	0.204	1.840	5.683	0.084
Zhou et al. [6]	No	70	0.383	5.321	10.47	0.478
Godard et al. [5]	No	70	0.544	10.94	11.76	0.193
<b>Ours</b>	<b>No</b>	<b>70</b>	<b>0.393</b>	<b>5.714</b>	<b>8.908</b>	<b>0.186</b>

Absolute Relative Difference (AbsRel), Squared Relative Difference (SqRel), Root Mean Square Error (RMSE) and Root Mean Squared Logarithmic Error (RMSElog). The accuracy metrics [4,23] that we use are the percentages of estimated depth  $d_{est}$  that subject to

$$\max\left(\frac{d_{est}}{d_{gt}}, \frac{d_{gt}}{d_{est}}\right) = \delta < threshold \quad (11)$$

**Results on KITTI dataset.** The evaluation results on the KITTI dataset are reported in Table 1. We use different combinations of train/test splits (E for Eigen, G for Godard) and cap distances (80m and 50m) to compare with different works. For Eigen et al. [13], Liu et al. [23], Zhou et al. [6] and Godard et al. [5], the Eigen split with 80m cap distance are used. For Garg et al. [4], Zhou et al. [6] and Godard et al. [5], the Eigen split with 50m cap distance are used. We also report our result on Godard split with 80m cap. For the ablation study of the ZNCC loss, we have implemented a patch-based Sum of Absolute Differences (SAD) loss that is a common and basic similarity measurement used for stereo matching algorithm to replace the ZNCC loss and keep the same multi-level patch setting. The results for the multi-level patch-based SAD

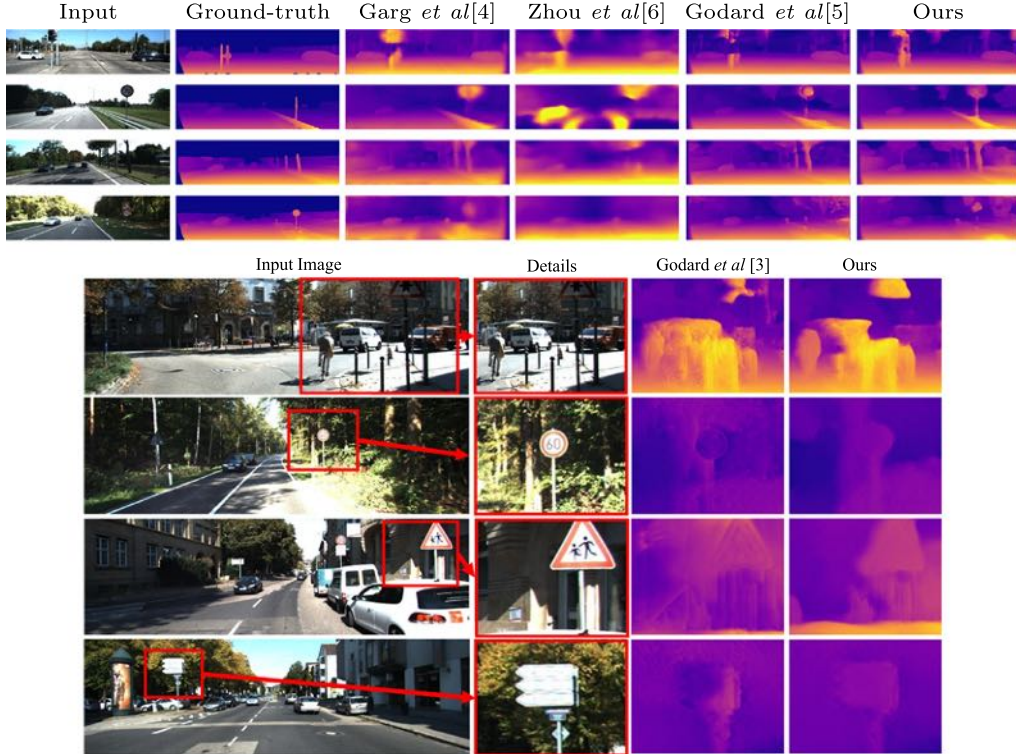
loss are reported as ours-SAD, which shows that our dedicated multi-level patch-based loss with SAD similarity measurement can already improve the benchmark results, but more improvements came with our proposed multi-level patch-based loss using the advanced ZNCC similarity measurement (reported as ours-ZNCC), which achieved the state-of-the-art results for monocular depth estimation problem on KITTI dataset.

**Results on Make3D dataset.** To further access the generalization ability of our proposed methods and compare with other methods, we also evaluate our trained networks on Make3D dataset [18]. For supervised methods [21,25,50], they are trained using ground truth depth data from the Make3D training set. For unsupervised methods [5,6] and ours, are trained on KITTI + Cityscapes datasets without the presence of any image from Make3D dataset. For evaluation, we measure the error metrics (AbsRel, SqRel, RMSE and RMSElog) using the test image with ground truth from Make3D dataset. As can be seen from Table 2, although our method scored similar results to Zhou et al. [6] regarding relative errors, for the RMSE, our methods outperform all of the state-of-the-art unsupervised methods.

Compared among unsupervised methods, our method produced better results regarding RMSE (RMSE and RMSElog) and at large cap distance (70m and 80m), and not significantly improve the relative error metrics (AbsRel, SqRel) at small cap distance (50m). This is totally what we expect as our multi-scale patch-based loss function performs better results when the distances of left-right corresponding pixels are large (meaning the pixel is at large distance), which the pixel-based loss function will prone to fail.

#### 4.3.2. Qualitative evaluation

The qualitative comparison to some of the related methods on KITTI dataset is shown in Fig. 6. While our network structure is



**Fig. 6.** Upper part: comparison of monocular depth estimation on KITTI dataset between Garg et al. [4], Zhou et al. [6], Godard et al. [5], and ours. Lower part: comparison of details with Godard et al. [5]. All of the results are generated using authors' provided pre-trained model. The ground-truth depth map is interpolated from sparse point map only for visualization.

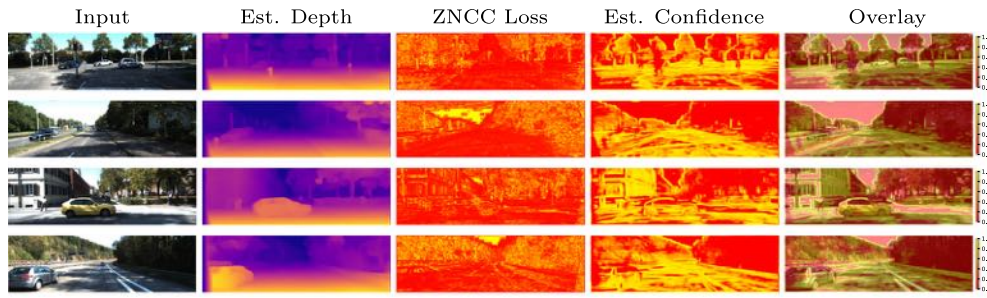


Fig. 7. Confidence estimation results. A colorbar from red to yellow is used to represent 0 to 1.

similar to that of Godard et al. [5], both generate clear and accurate depth than other works. We also provide a detailed comparison with the results of Godard et al. [5] in the lower part of Fig. 6. Our network can generate more accurate depth in complex regions with thin structures and texture-less areas such as the pillars and traffic signs. This verified the theory we explained in Fig. 5 that our patch-based loss function is more robust and easier to converge to the global minimum in complex regions.

#### 4.3.3. Confidence map evaluation

We show the confidence estimation results in Fig. 7. A colorbar from red to yellow is used to represent 0 to 1. We can see that the estimated confidence can nicely represent the inverted ZNCC loss but less noisy due to the small network we use to prevent over-fitting. The overlaid confidence on input image shows that our ConfidenceNet has learned to generate confidence from contextual information. For example, in texture-less areas (sky, building), dark areas (trees under shadow), occluded areas (around thin structures) and reflective areas (car window), the estimated confidence is usually very low, while the texture-rich areas and edges usually have high confidence.

## 5. Discussion

In this paper, we have presented a novel self-supervised framework for monocular depth learning and confidence estimation. We incorporate the patch matching theory into a fully differential DCNN and achieve self-supervised training of both depth and the confidence of depth. Our proposed loss function exploits the epipolar constraint of stereo vision and also provides a normalized similarity that is further used to supervise the confidence estimation. Our method not only outperforms the state-of-the-art results on the KITTI benchmark evaluation, but also for the first time, we are able to simultaneously generate depth from monocular images and estimate the confidence of the generated depth. This is a step change for monocular depth estimation as it significantly increases the feasibility of using monocular depth estimation into many practical applications such as autonomous driving and monocular endoscopic surgery [7], where the accuracy of estimated depth is crucial.

**Why Our ConfidenceNet works?** As there is certain limitation of unsupervised monocular depth learning from stereo pairs (ambiguous depth estimation in texture-less area, reflection, etc.). Our ConfidenceNet is supervised by the per-pixel ZNCC loss of our depth estimation network (which can be regarded as the confidence of current depth), it explicitly learns the regions where our depth estimation network performs well and badly. But on a deeper level, our ConfidenceNet actually implicitly learns the inherent defect of the patch matching algorithm – it would fail on texture-less regions and performs badly near stereo view occlusions, reflections and blurred areas. Therefore, after sufficient train-

ing steps, our ConfidenceNet can capture and memory where the DepthNet would perform good or bad, and give an estimation of the confidence of our DepthNet, although they are two different networks.

**In future work.** We will continue optimizing our model and explore the possibility of using adaptive window size for patch sampling to decrease the training time and increase accuracy in small structures.

## Acknowledgments

This research was supported/partially supported by the EU H2020 Programme, H2020-MSCA-RISE-2018:iGame project 82387.

## References

- [1] B.J. Dunkin, C. Flowers, 3D in the minimally invasive surgery (MIS) operating room: cameras and displays in the evolution of MIS, in: *Imaging and Visualization in The Modern Operating Room*, Springer, 2015, pp. 145–155.
- [2] A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks, in: F. Pereira, C.J.C. Burges, L. Bottou, K.Q. Weinberger (Eds.), *Advances in Neural Information Processing Systems 25*, Curran Associates, Inc., 2012, pp. 1097–1105.
- [3] Z. Zhang, Determining the epipolar geometry and its uncertainty: a review, *Int. J. Comput. Vis.* 27 (2) (1998) 161–195, doi:10.1023/A:1007941100561.
- [4] R. Garg, V.K. B.G., G. Carneiro, I. Reid, Unsupervised CNN for single view depth estimation: geometry to the rescue, in: B. Leibe, J. Matas, N. Sebe, M. Welling (Eds.), *Proceedings of the Computer Vision – ECCV, Springer International Publishing, Cham*, 2016, pp. 740–756.
- [5] C. Godard, O.M. Aodha, G.J. Brostow, Unsupervised monocular depth estimation with left-right consistency, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6602–6611, doi:10.1109/CVPR.2017.699.
- [6] T. Zhou, M. Brown, N. Snavely, D.G. Lowe, Unsupervised learning of depth and ego-motion from video, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6612–6619, doi:10.1109/CVPR.2017.700.
- [7] L. Chen, W. Tang, N.W. John, T.R. Wan, J.J. Zhang, Slam-based dense surface reconstruction in monocular minimally invasive surgery and its application to augmented reality, *Comput. Methods Progr. Biomed.* 158 (2018) 135–146, doi:10.1016/j.cmpb.2018.02.006.
- [8] K. Tateno, F. Tombari, I. Laina, N. Navab, Cnn-slam: real-time dense monocular slam with learned depth prediction, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6565–6574, doi:10.1109/CVPR.2017.695.
- [9] S.T. Barnard, M.A. Fischler, Computational stereo, *ACM Comput. Surv.* 14 (4) (1982) 553–572, doi:10.1145/356893.356896.
- [10] D. Scharstein, R. Szeliski, R. Zabih, A taxonomy and evaluation of dense two-frame stereo correspondence algorithms, in: *Proceedings IEEE Workshop on Stereo and Multi-Baseline Vision (SMBV)*, 2001, pp. 131–140, doi:10.1109/SMBV.2001.988771.



- [11] H. Hirschmüller, Stereo processing by semiglobal matching and mutual information, *IEEE Trans. Pattern Anal. Mach. Intelligence* 30 (2) (2008) 328–341, doi:10.1109/TPAMI.2007.1166.
- [12] A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, End-to-end learning of geometry and context for deep stereo regression, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 66–75, doi:10.1109/ICCV.2017.17.
- [13] D. Eigen, C. Puhrsch, R. Fergus, Depth map prediction from a single image using a multi-scale deep network, in: *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, in: *NIPS'14*, MIT Press, Cambridge, MA, USA, 2014, pp. 2366–2374.
- [14] R. Zhang, P.-S. Tsai, J.E. Cryer, M. Shah, Shape-from-shading: a survey, *IEEE Trans. Pattern Anal. Mach. Intell.* 21 (8) (1999) 690–706, doi:10.1109/34.784284.
- [15] A. Saxena, S.H. Chung, A.Y. Ng, Learning depth from single monocular images, in: Y. Weiss, B. Schölkopf, J.C. Platt (Eds.), *Advances in Neural Information Processing Systems 18*, MIT Press, 2006, pp. 1161–1168.
- [16] A. Saxena, J. Schulte, A.Y. Ng, Depth estimation using monocular and stereo cues, in: *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, in: *IJCAI'07*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2007, pp. 2197–2203.
- [17] A. Saxena, S.H. Chung, A.Y. Ng, 3-D depth reconstruction from a single still image, *Int. J. Comput. Vis.* 76 (1) (2008) 53–69, doi:10.1007/s11263-007-0071-y.
- [18] A. Saxena, M. Sun, A.Y. Ng, Make3d: learning 3D scene structure from a single still image, *IEEE Trans. Pattern Anal. Mach. Intell.* 31 (5) (2009) 824–840, doi:10.1109/TPAMI.2008.132.
- [19] D. Eigen, R. Fergus, Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 2650–2658, doi:10.1109/ICCV.2015.304.
- [20] B. Li, C. Shen, Y. Dai, A. van den Hengel, M. He, Depth and surface normal estimation from monocular images using regression on deep features and hierarchical crfs, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1119–1127, doi:10.1109/CVPR.2015.7298715.
- [21] M. Liu, M. Salzmann, X. He, Discrete-continuous depth estimation from a single image, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 716–723, doi:10.1109/CVPR.2014.97.
- [22] Y. Hua, H. Tian, Depth estimation with convolutional conditional random field network, *Neurocomputing* 214 (2016) 546–554, doi:10.1016/j.neucom.2016.06.029.
- [23] F. Liu, C. Shen, G. Lin, I. Reid, Learning depth from single monocular images using deep convolutional neural fields, *IEEE Trans. Pattern Anal. Mach. Intell.* 38 (10) (2016) 2024–2039, doi:10.1109/TPAMI.2015.2505283.
- [24] D. Xu, E. Ricci, W. Ouyang, X. Wang, N. Sebe, Multi-scale continuous CRFS as sequential deep networks for monocular depth estimation, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 161–169, doi:10.1109/CVPR.2017.25.
- [25] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, N. Navab, Deeper depth prediction with fully convolutional residual networks, in: *Proceedings of the Fourth International Conference on 3D Vision (3DV)*, 2016, pp. 239–248, doi:10.1109/3DV.2016.32.
- [26] Y. Cao, Z. Wu, C. Shen, Estimating depth from monocular images as classification using deep fully convolutional residual networks, *IEEE Trans. Circuits Syst. Video Technol.* PP (99) (2017) 1, doi:10.1109/TCSVT.2017.2740321.
- [27] J. Li, R. Klein, A. Yao, A two-streamed network for estimating fine-scaled depth maps from single RGB images, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 3392–3400, doi:10.1109/ICCV.2017.365.
- [28] L. Ladický, J. Shi, M. Pollefeys, Pulling things out of perspective, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, in: *CVPR '14*, IEEE Computer Society, Washington, DC, USA, 2014, pp. 89–96, doi:10.1109/CVPR.2014.19.
- [29] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, A. Yuille, Towards unified depth and semantic prediction from a single image, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 2800–2809, doi:10.1109/CVPR.2015.7298897.
- [30] A. Mousavian, H. Pirsiaavash, J. Košecká, Joint semantic segmentation and depth estimation with deep convolutional networks, in: *Proceedings of the Fourth International Conference on 3D Vision (3DV)*, IEEE, 2016, pp. 611–619.
- [31] H. Yan, S. Zhang, Y. Zhang, L. Zhang, Monocular depth estimation with guidance of surface normal map, *Neurocomputing* 280 (2018) 86–100, doi:10.1016/j.neucom.2017.08.074.
- [32] Y. Kuznetsov, J. Stckler, V. Leibe, Semi-supervised deep learning for monocular depth map prediction, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2215–2223, doi:10.1109/CVPR.2017.238.
- [33] F. Ma, S. Karaman, Sparse-to-dense: depth prediction from sparse depth samples and a single image, in: *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 1–8, doi:10.1109/ICRA.2018.8460184.
- [34] D. Zoran, P. Isola, D. Krishnan, W.T. Freeman, Learning ordinal relationships for mid-level vision, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 388–396, doi:10.1109/ICCV.2015.52.
- [35] W. Chen, Z. Fu, D. Yang, J. Deng, Single-image depth perception in the wild, in: D.D. Lee, M. Sugiyama, U.V. Luxburg, I. Guyon, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 29*, Curran Associates, Inc., 2016, pp. 730–738.
- [36] A. Fitzgibbon, Y. Wexler, A. Zisserman, Image-based rendering using image-based priors, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2003, pp. 1176–1183vol.2, doi:10.1109/ICCV.2003.1238625.
- [37] T. Zhou, S. Tulsiani, W. Sun, J. Malik, A.A. Efros, View synthesis by appearance flow, in: *Proceedings of the European Conference on Computer Vision*, 2016.
- [38] J. Flynn, I. Neulander, J. Philbin, N. Snavely, Deep stereo: Learning to predict new views from the world's imagery, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 5515–5524, doi:10.1109/CVPR.2016.595.
- [39] J. Xie, R. Girshick, A. Farhadi, Deep3D: dully automatic 2D-to-3D video conversion with deep convolutional neural networks, in: B. Leibe, J. Matas, N. Sebe, M. Welling (Eds.), *Proceedings of the Computer Vision – ECCV*, Springer International Publishing, Cham, 2016, pp. 842–857.
- [40] C. Wang, J. Miguel Buenaposada, R. Zhu, S. Lucey, Learning depth from monocular videos using direct methods, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [41] R. Li, S. Wang, Z. Long, D. Gu, Undeepvo: monocular visual odometry through unsupervised deep learning, in: *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2018, pp. 7286–7291.
- [42] H. Zhan, R. Garg, C. Saroj Weerasekera, K. Li, H. Agarwal, I. Reid, Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 340–349.
- [43] Z. Li, N. Snavely, Megadepth: Learning single-view depth prediction from internet photos, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2041–2050.
- [44] X. Guo, H. Li, S. Yi, J. Ren, X. Wang, Learning monocular depth by distilling cross-domain stereo networks, in: V. Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), *Proceedings of the Computer Vision – ECCV*, Springer International Publishing, Cham, 2018, pp. 506–523.
- [45] Y. Luo, J. Ren, M. Lin, J. Pang, W. Sun, H. Li, L. Lin, Single view stereo matching, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [46] A. Dosovitskiy, J.T. Springenberg, M. Tatarchenko, T. Brox, Learning to generate chairs, tables and cars with convolutional networks, *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (4) (2017) 692–705, doi:10.1109/TPAMI.2016.2567384.
- [47] N. Mayer, E. Ilg, P. Husser, P. Fischer, D. Cremers, A. Dosovitskiy, T. Brox, A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4040–4048, doi:10.1109/CVPR.2016.438.
- [48] E. Shelhamer, J. Long, T. Darrell, Fully convolutional networks for semantic segmentation, *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (4) (2017) 640–651, doi:10.1109/TPAMI.2016.2572683.
- [49] M. Jaderberg, K. Simonyan, A. Zisserman, k. kavukcuoglu, Spatial transformer networks, in: C. Cortes, N.D. Lawrence, D.D. Lee, M. Sugiyama, R. Garnett (Eds.), *Advances in Neural Information Processing Systems 28*, Curran Associates, Inc., 2015, pp. 2017–2025.
- [50] K. Karsch, C. Liu, S.B. Kang, Depth transfer: depth extraction from video using non-parametric sampling, *IEEE Trans. Pattern Anal. Mach. Intell.* 36 (11) (2014) 2144–2158, doi:10.1109/TPAMI.2014.2316835.